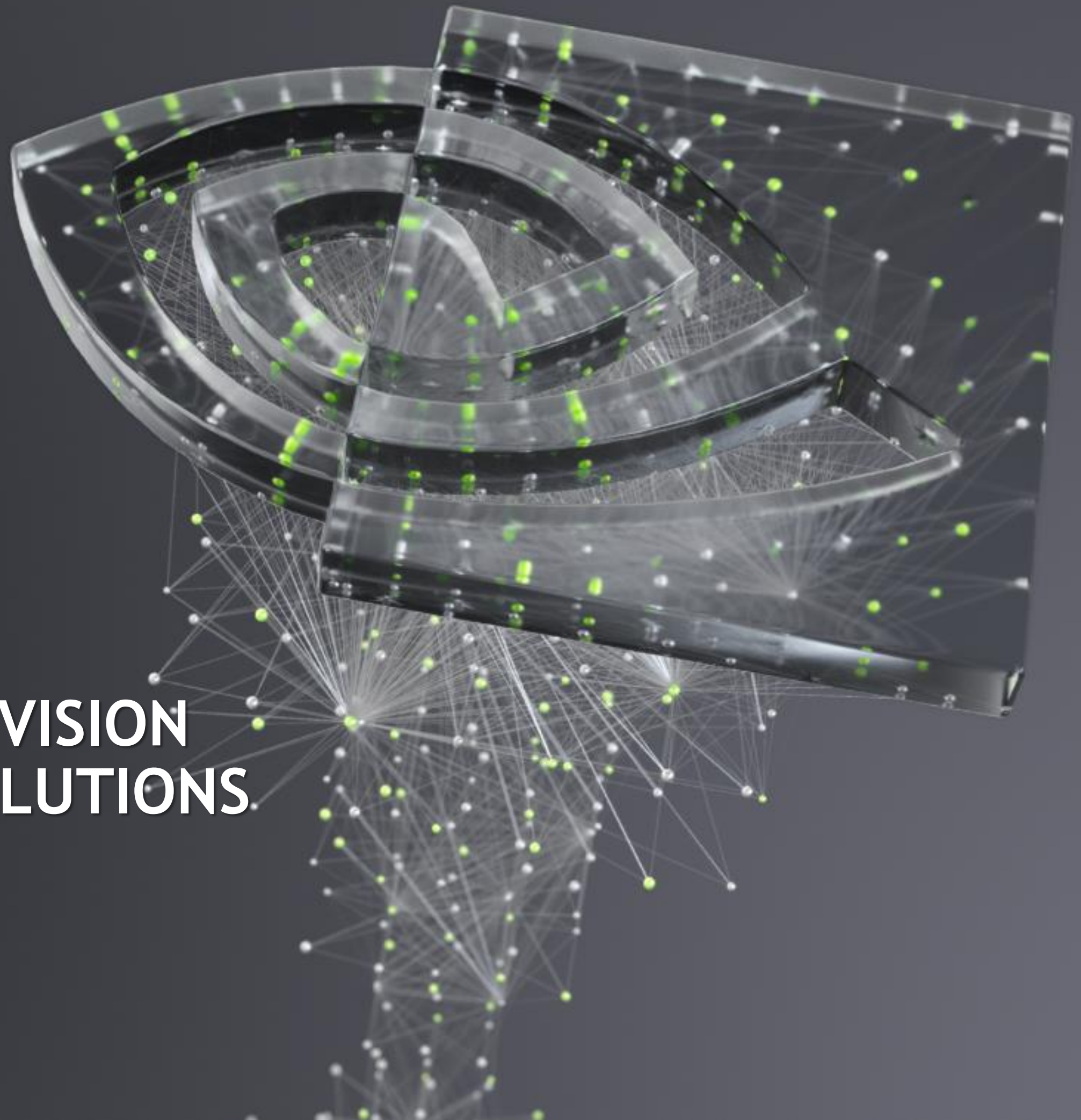




IMPLEMENTING COMPUTER VISION AND IMAGE PROCESSING SOLUTIONS WITH VPI

Rodolfo Lima, Feb. 11th, 2021



AGENDA

What we will cover today

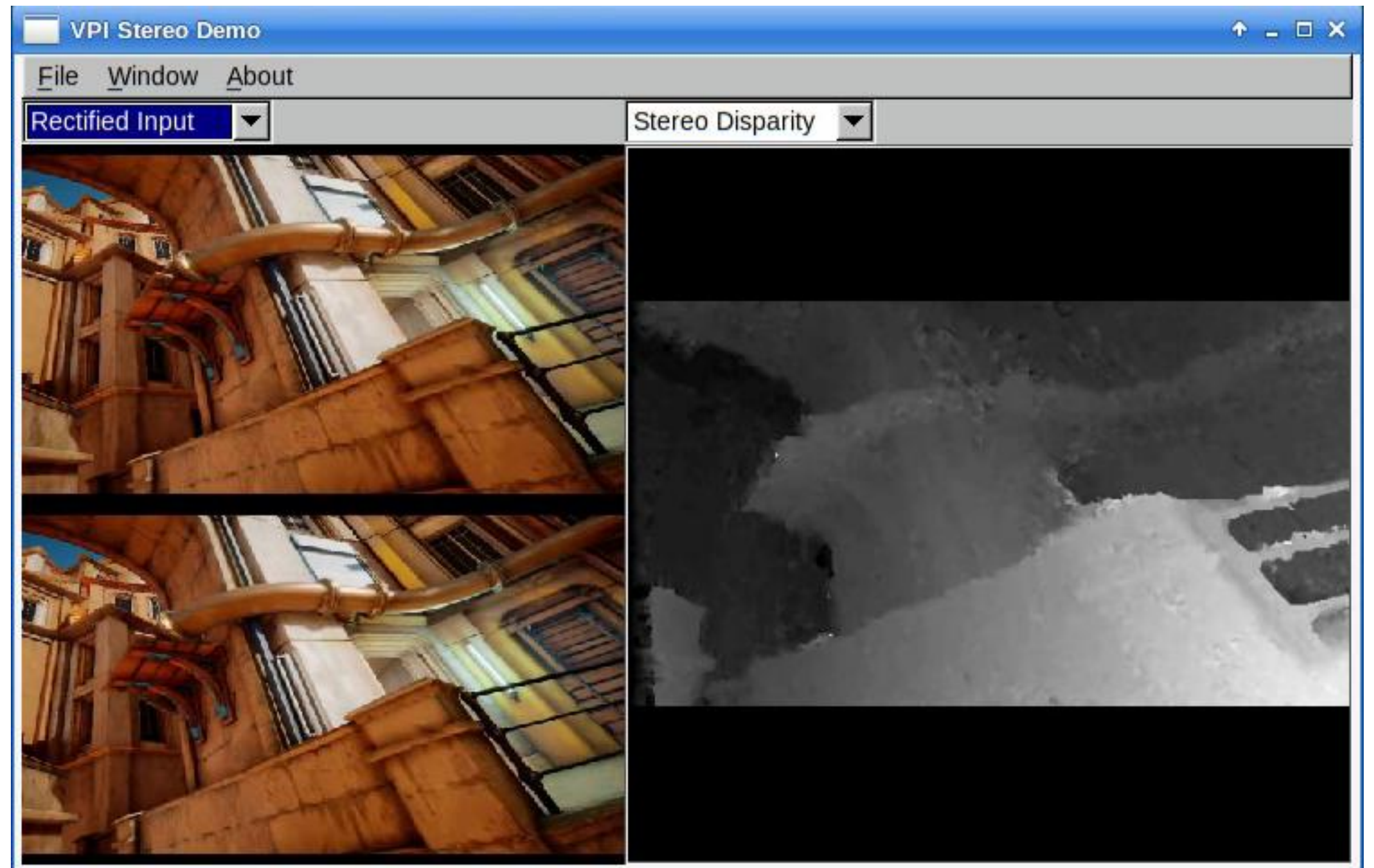
Why VPI?

Core Concepts

Example: Stereo Disparity Estimation

Future Work

Q & A



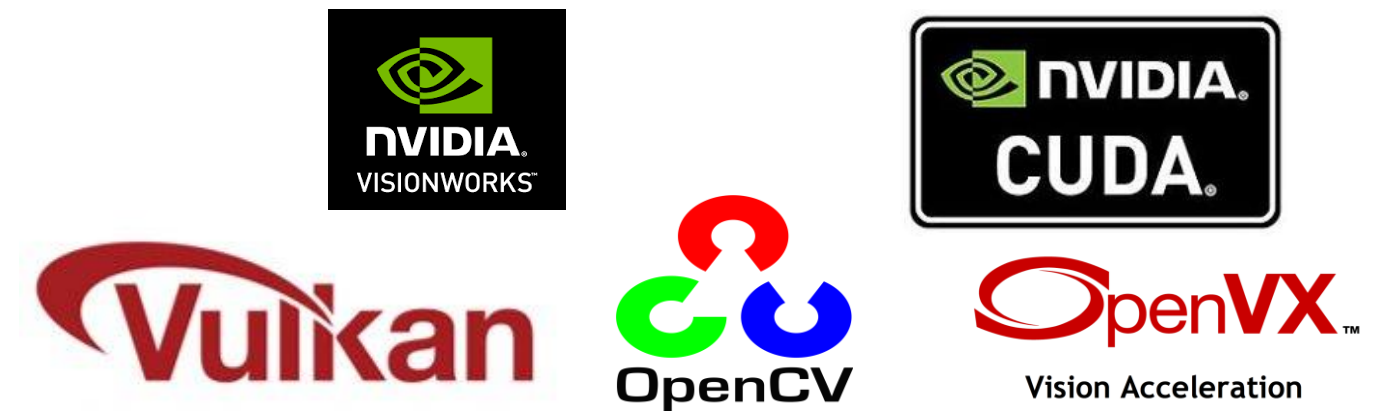


WHY VPI?

PROBLEM STATEMENT

What makes current solutions too complicated?

- ▶ Multiple mutually **incompatible** APIs needed
- ▶ Efficient memory management **hard** to get right
- ▶ **Difficult** to experiment with the code or fine-tune it
- ▶ Some compute engines **aren't directly accessible** via a public API.



VISION PROGRAMMING INTERFACE - VPI

NVIDIA's next-gen API for embedded CV

- ▶ Create **efficient** CV pipelines with all Jetson embedded accelerators
- ▶ Same algorithm implemented by different accelerators
- ▶ Easily **load-balance** CV workloads at system level
- ▶ **Unified API** to interface with different accelerators
- ▶ Accelerate on both **Jetson and x86 Linux PCs**
- ▶ **Zero-copy** memory management among different accelerators*
- ▶ **Interoperability** with OpenCV, NVIDIA® CUDA®, EGL, among others
- ▶ Designed to **supersede** NVIDIA® VisionWorks™

* When allowed by backend and memory characteristics



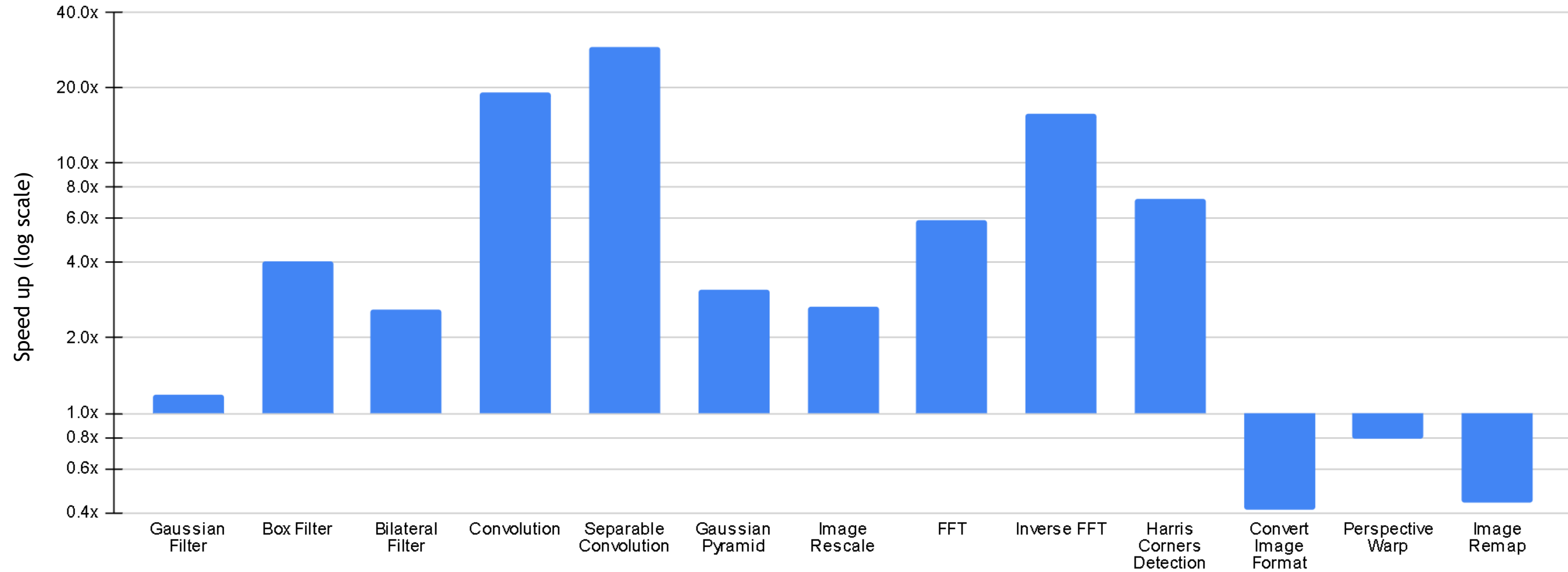
The screenshot displays the NVIDIA VPI website interface. At the top left is the NVIDIA logo. The page title is "VPI - Vision Programming Interface" with "1.0 Release" on the right. A navigation bar includes "Main Page", "Related Pages", "API Reference", and a search box. A left sidebar lists the following items: "VPI - Vision Programming Interface", "Release Notes v1.0", "Getting Started", "Architecture", "Performance Benchmark", "Algorithms" (with sub-items: Box Filter, Bilateral Filter, Gaussian Filter, Gaussian Pyramid Generator, Convolution, Separable Convolution, Convert Image Format, Rescale, Remap, Perspective Warp, FFT, Inverse FFT, Lens Distortion Correction, Stereo Disparity Estimator, KLT Feature Tracker, Harris Corner Detector, Temporal Noise Reduction, Pyramidal LK Optical Flow), "Sample Applications", "Appendices", and "End User License Agreement". The main content area features the text "replacement of existing computing tasks with faster VPI equivalents." and "Here are some examples of algorithms provided by VPI:". Below this are four image examples: "Stereo Disparity Estimator" (a grayscale image with a depth map overlay), "KLT Feature Tracker" (a street scene with tracked features), "Harris Corner Detector" (a grayscale image of a building facade with red corner markers), and "Remap" (a 360-degree panoramic view of a globe).

CPU PERFORMANCE BENCHMARK

Significant Speed Up Compared to OpenCV

Up to 30x faster than OpenCV on CPU

Jetson AGX Xavier

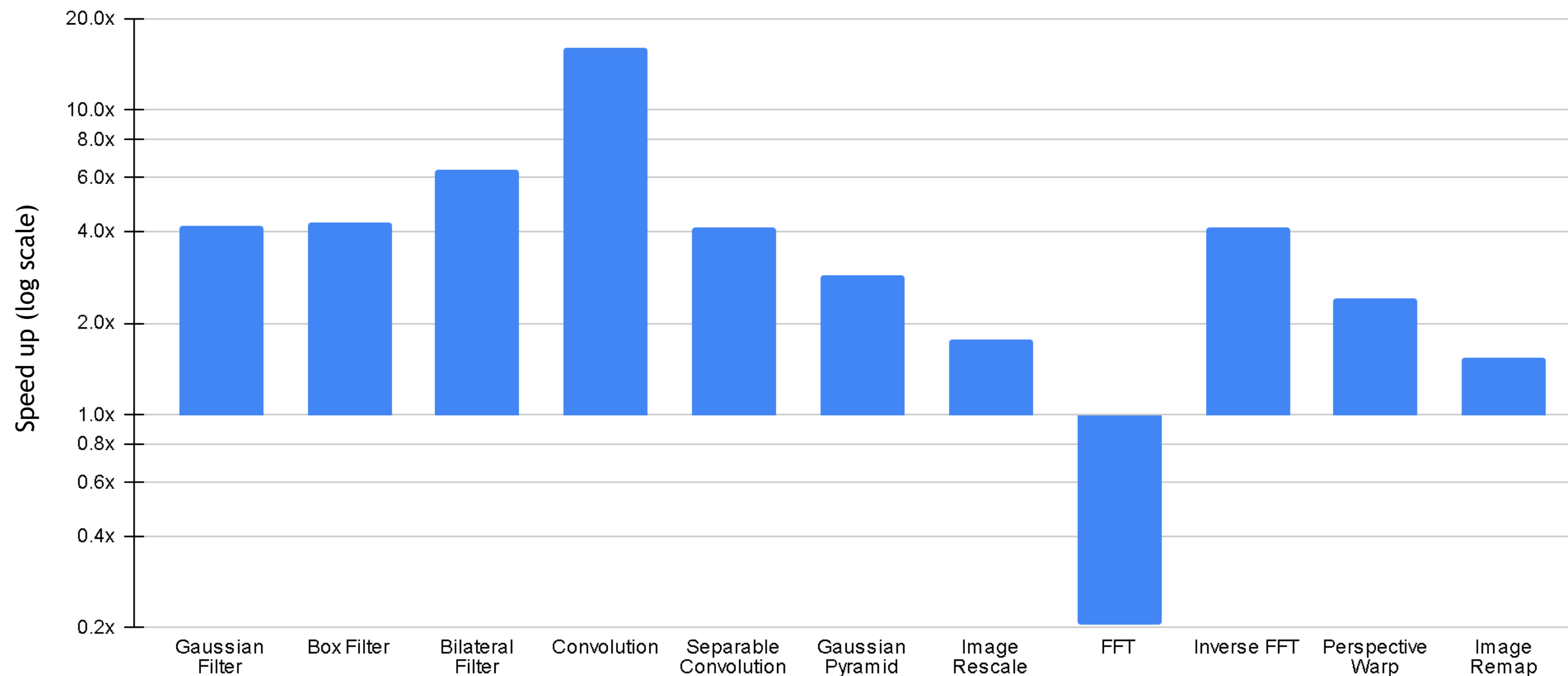


GPU PERFORMANCE BENCHMARK

Significant Speed Up Compared to OpenCV

Up to 15x faster than OpenCV on GPU

Jetson AGX Xavier

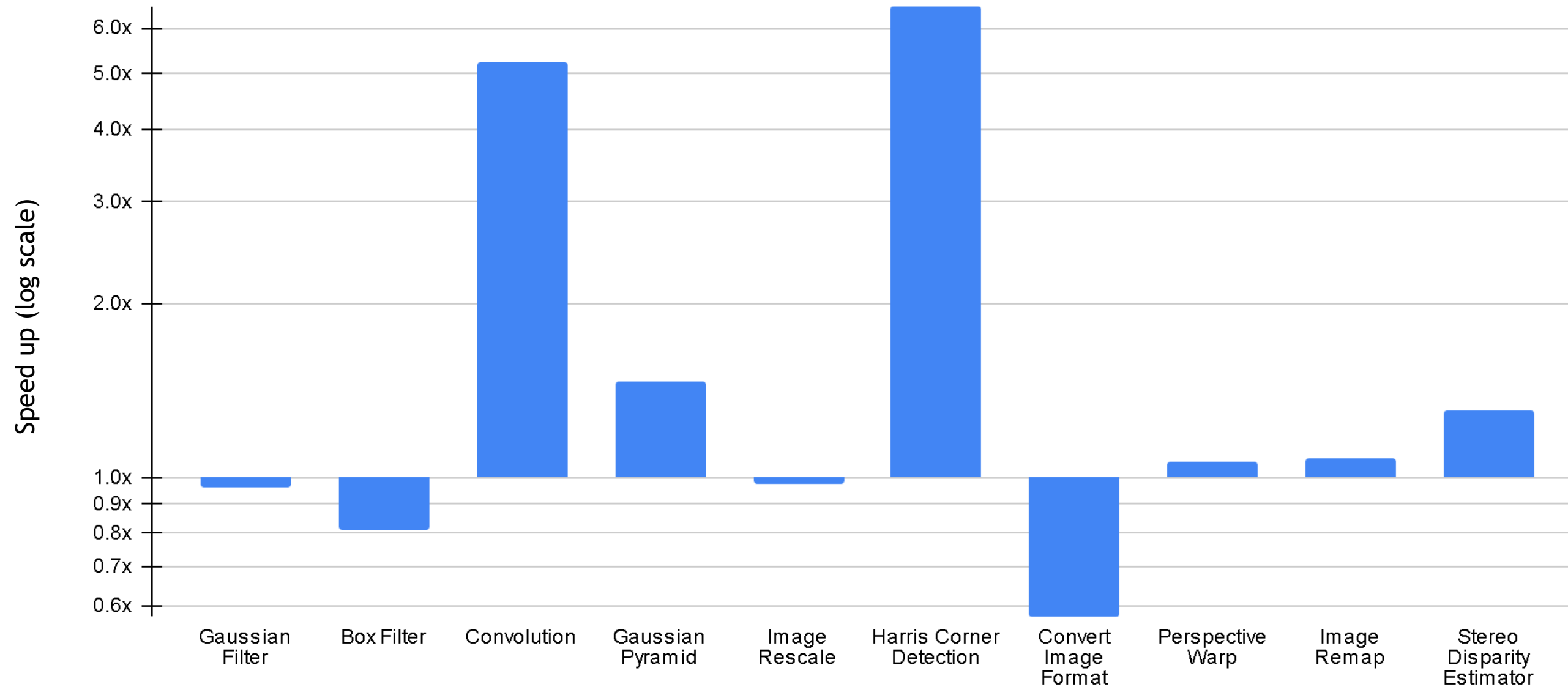


GPU PERFORMANCE BENCHMARK

Significant Speed Up Compared to VisionWorks

Up to 6x faster than VisionWorks

Jetson AGX Xavier

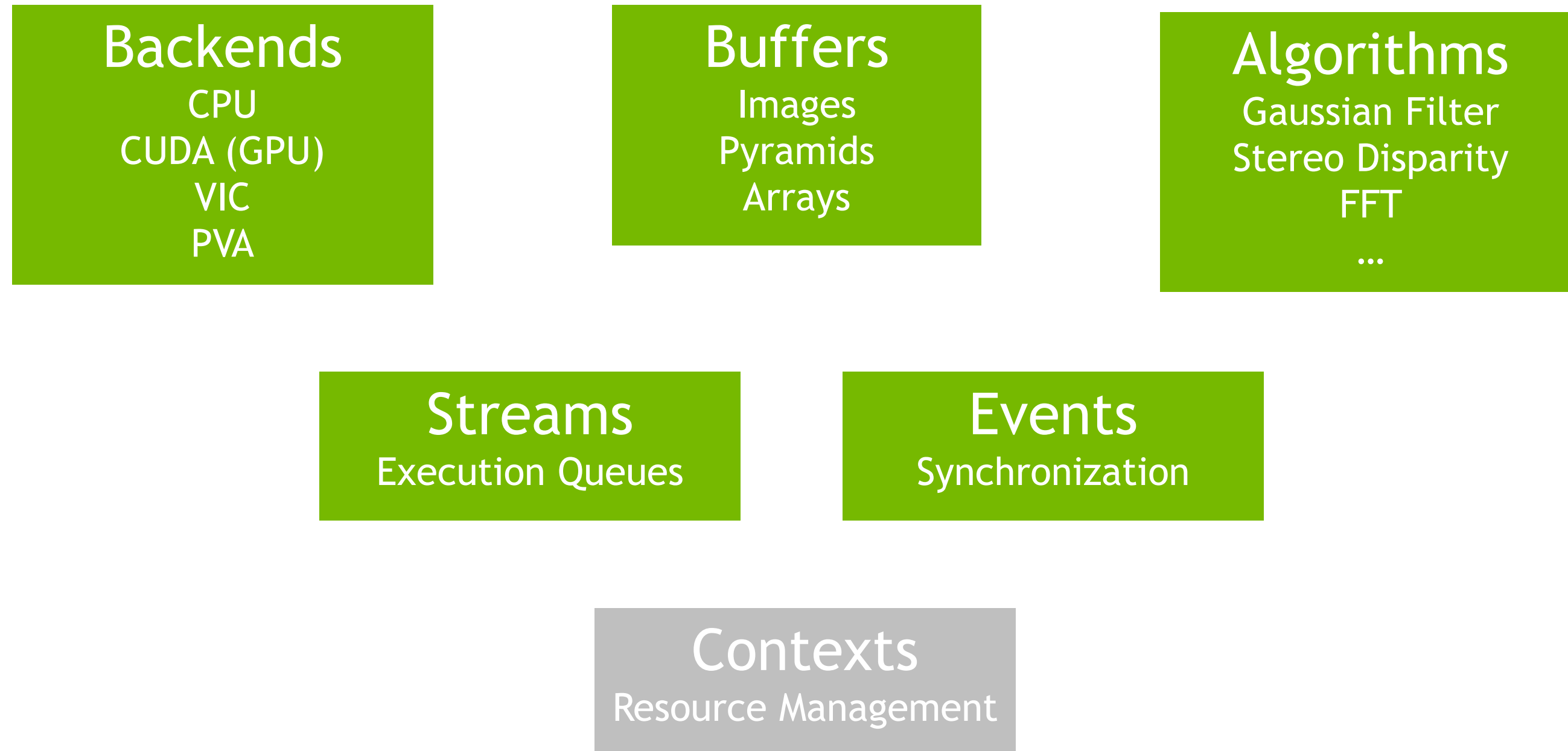




CORE CONCEPTS

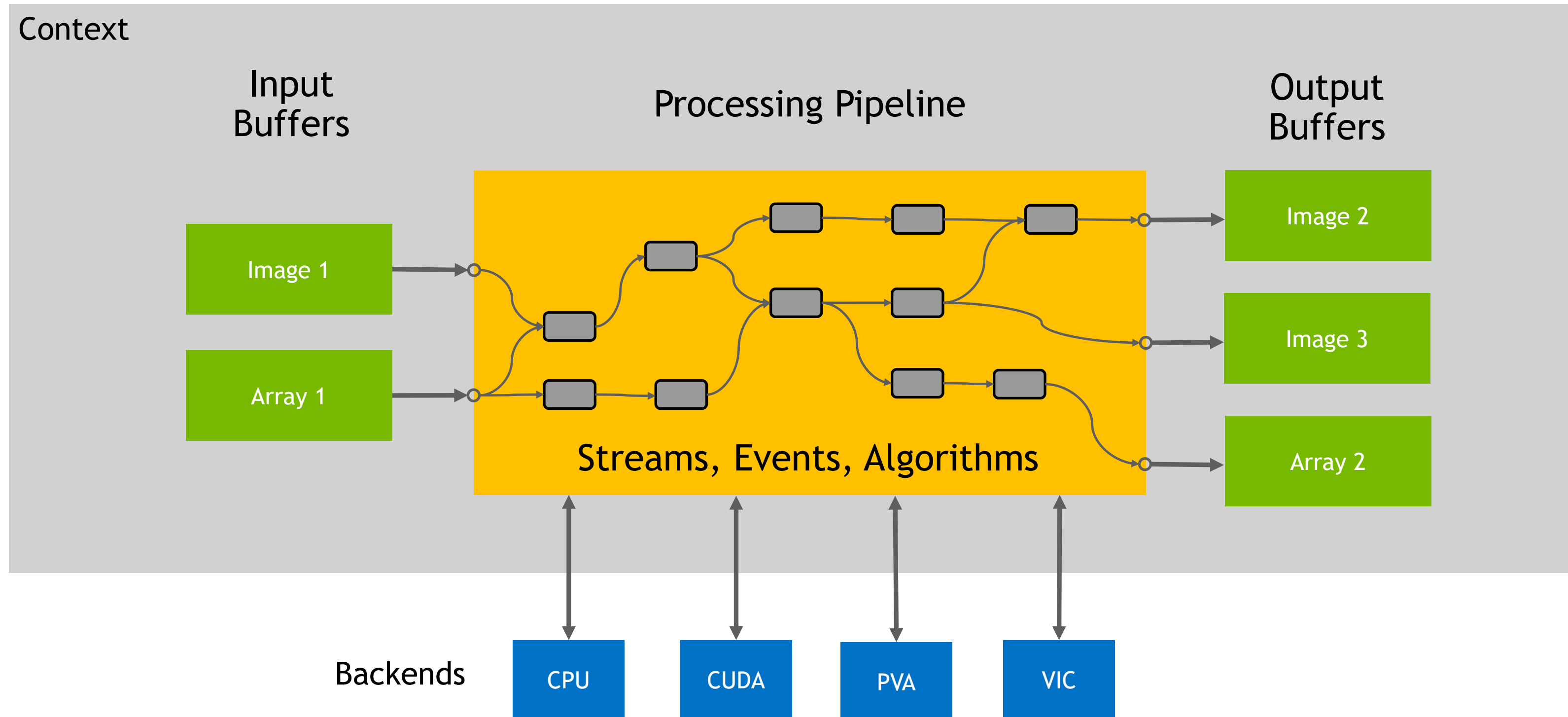
BUILDING BLOCKS

Main VPI components



PROCESSING PIPELINE

How components fit together



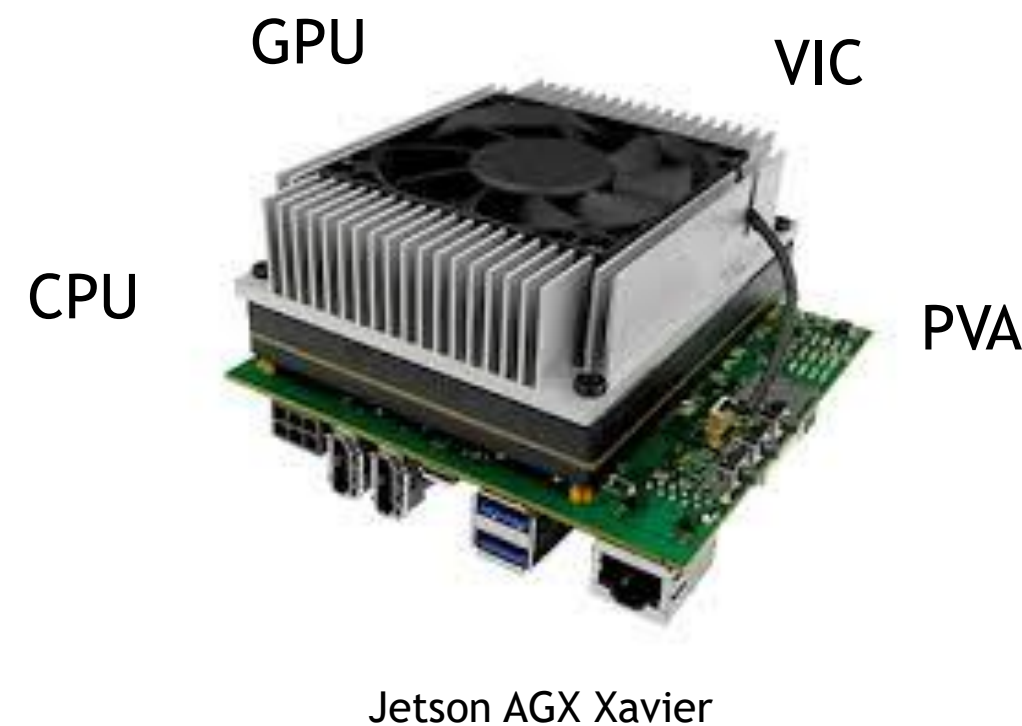
BACKENDS

Overview

Backends

CPU
CUDA
PVA
VIC

- ▶ A backend represents one type of compute engine
- ▶ Certain tasks are more fit to some specific backend
- ▶ Same algorithm implemented in multiple backends
 - ▶ Allows offloading tasks to less used engines
 - ▶ Eg: leave GPU free for DL inference



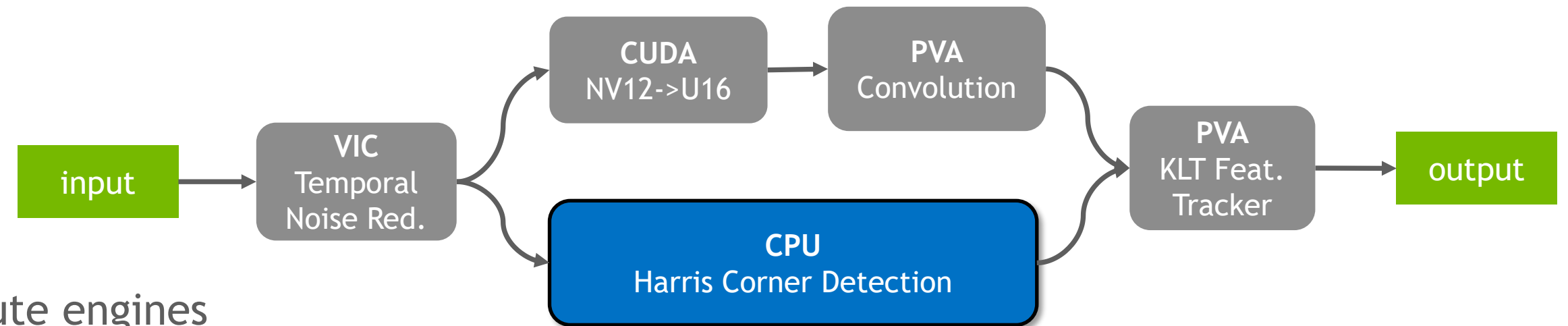
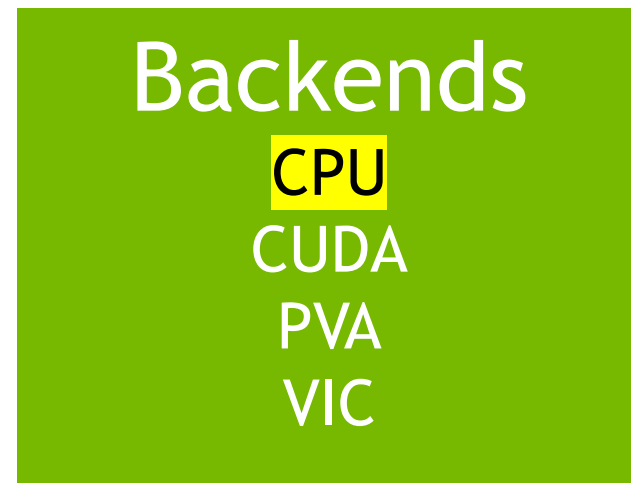
Jetson AGX Xavier



Geforce RTX 2070 + CPU

BACKENDS

CPU - Central Processing Unit



- ▶ Offload computation from other compute engines
- ▶ CPU would be left idle otherwise
- ▶ More cores -> faster algorithm execution
- ▶ Available on Jetson/aarch64 and x86 architectures

BACKENDS

CUDA - Compute Unified Device Architecture

Backends

CPU
CUDA
PVA
VIC

- ▶ Represents the GPU on the system
- ▶ Speed
- ▶ Flexibility - less restrictions on algorithm parameters
- ▶ Available on Jetson and x86 architectures
- ▶ Supports Maxwell, Pascal, Volta, Turing and Ampere GPU architectures



Ampere GA102 GPU HW architecture

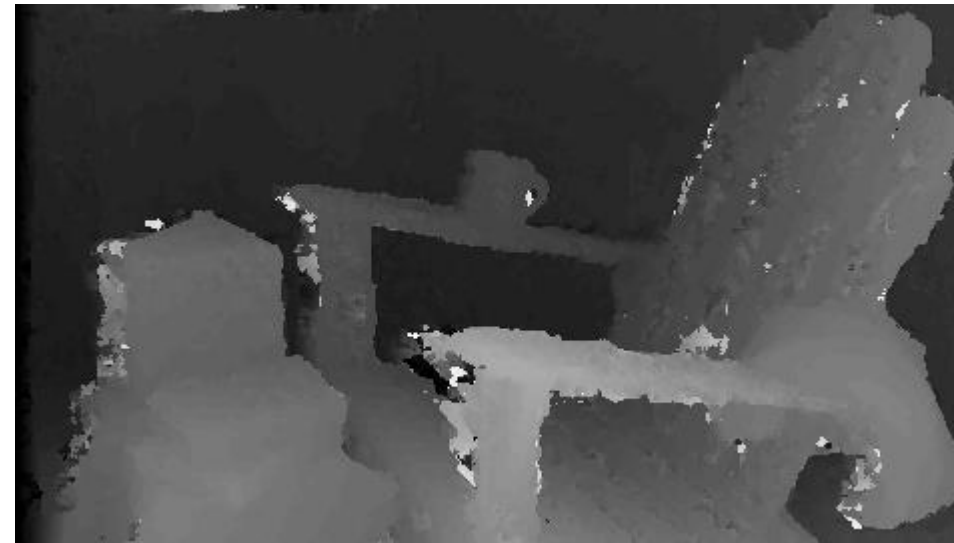
BACKENDS

PVA - Programmable Vision Accelerator

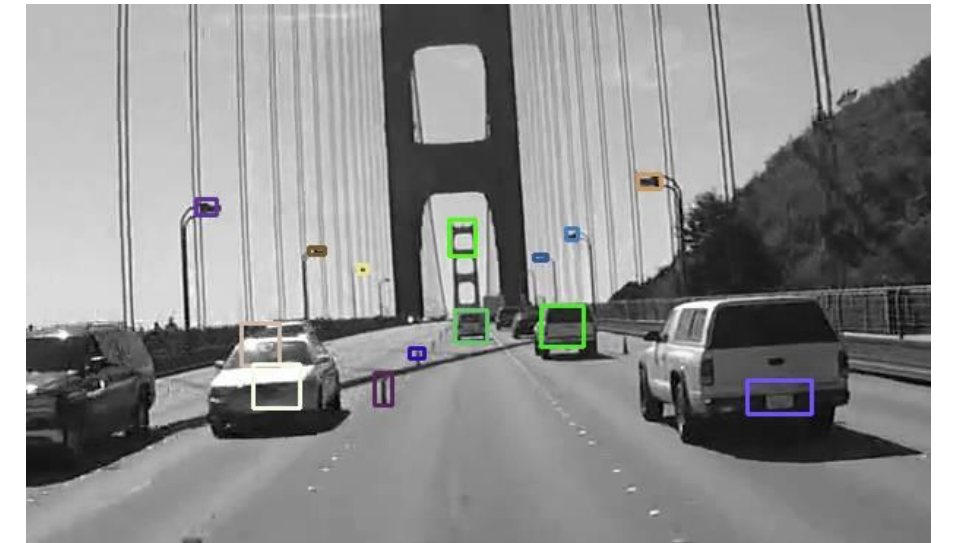
Backends

CPU
CUDA
PVA
VIC

- ▶ Used when power efficiency is required
- ▶ Available on Jetson AGX Xavier and Jetson Xavier NX devices.
- ▶ Four independent processors per device
- ▶ Increased task parallelism
- ▶ Good choice to offload processing from GPUs.



Stereo Disparity Estimator



KLT Feature Tracker



Harris Corner Detection



Gaussian Filter

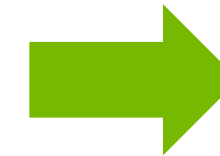
BACKENDS

VIC - Video Image Compositor

Backends

CPU
CUDA
PVA
VIC

- ▶ Fixed-function processor
- ▶ Low power usage
- ▶ Optimized for video processing
- ▶ Available on all Jetson devices
- ▶ Some operations available only on Jetson AGX Xavier and Jetson Xavier NX devices



Perspective Warp



Remap

BUFFERS

VPIImage overview

Buffers

Image

Pyramids

Arrays

- ▶ Images are characterized by:
 - ▶ Dimensions (width x height)
 - ▶ Image Format - NV12, U8, RGB8, ...
 - ▶ Backends that can use them: CUDA, CPU, PVA, ...



BUFFERS

VPIImage usage

Buffers

Image

Pyramids

Arrays

- ▶ Images are characterized by:
 - ▶ Dimensions (width x height)
 - ▶ Image Format - NV12, U8, RGB8, ...
 - ▶ Backends that can use them: CUDA, CPU, PVA, ...

```
/* Create a 512x256 RGBA8 VPIImage
   to be used by CPU and CUDA backends. */
1. VPIImage img;
2. vpiImageCreate (512, 256,
3.                 VPI_IMAGE_FORMAT_RGBA8,
4.                 VPI_BACKEND_CPU | VPI_BACKEND_CUDA,
5.                 &img);

6. /* (omitted) use img */

/* When finished, destroy/deallocate it */
7. vpiImageDestroy (img);
```

BUFFERS

VPIImage - OpenCV Interoperability

Buffers

Image

Pyramids

Arrays

- ▶ Images can wrap user-allocated memory buffers:
 - ▶ **OpenCV**
 - ▶ CUDA device pointer (cudaMalloc, ...)
 - ▶ Raw CPU pointer (malloc, new, ...)
 - ▶ EGLImage*
 - ▶ NvBuffer*

```
/* Load an image from disk using OpenCV */
1. cv::Mat cvImage = cv::imread("baboon.png");

/* Wrap it to be used by VPI.
   Image format inferred by VPI = BGR8 */
2. VPIImage img;
3. vpiImageCreateOpenCVMatWrapper (cvImage,
4.     VPI_BACKEND_CPU|VPI_BACKEND_CUDA,
5.     &img);

/* Set wrapper to another cv::Mat */
6. vpiImageSetWrappedOpenCVMat (img, cvImage2);

/* When finished, destroy the wrapper */
7. vpiImageDestroy (img);

/* Now cvImage can be deallocated */
8. cvImage = cv::Mat();
```

* On Jetson devices only

BUFFERS

VPIImage - CUDA Interoperability

Buffers

Image

Pyramids

Arrays

- ▶ Images can wrap user-allocated memory buffers:
 - ▶ OpenCV
 - ▶ CUDA device pointer (cudaMalloc, ...)
 - ▶ Raw CPU pointer (malloc, new, ...)
 - ▶ EGLImage*
 - ▶ NvBuffer*

```
    /* Allocate a 2D CUDA memory buffer for RGBA8 data */
1. void *devBuffer;
2. size_t pitch;
3. cudaMallocPitch(&devBuffer, &pitch, 512*4, 256);

    /* Wrap it to be used by VPI. */
4. VPIImageData data;
5. memset(&data, 0, sizeof(data));
6. data.format = VPI_IMAGE_FORMAT_RGBA8;
7. data.numPlanes = 1;
8. data.planes[0].data = devBuffer;
9. data.planes[0].width = 512;
10. data.planes[0].height = 256;
11. data.planes[0].pitchBytes = pitch;
12. VPIImage img;
13. vpiImageCreateCUDAMemWrapper(&data,
14.                               VPI_BACKEND_CPU|VPI_BACKEND_CUDA,
15.                               &img);

    /* Set wrapper to a another CUDA memory buffer */
16. vpiImageSetWrappedCUDAMem(img, &data2);

    /* When finished with it, destroy it */
17. vpiImageDestroy(img);

    /* Now devBuffer can be deallocated */
18. cudaFree(devBuffer);
```

* On Jetson devices only

BUFFERS

VPIPyramid

Buffers
Image
Pyramids
Arrays

- ▶ Pyramids are characterized by:
 - ▶ Finer level dimensions (width x height)
 - ▶ Image Format - NV12, U8, RGB8, ...
 - ▶ Number of levels
 - ▶ Scale between one level and the next
 - ▶ Backends that can use them: CUDA, CPU, PVA, ...
 - ▶ More info: https://docs.nvidia.com/vpi/architecture.html#arch_pyramid



4-level Image Pyramid

BUFFERS

VPIArray

Buffers

Image
Pyramids
Arrays

- ▶ Arrays are characterized by:
 - ▶ Capacity (maximum number of elements)
 - ▶ Array type - key points, scores, bounding boxes ...
 - ▶ Number of elements (size), which can be redefined after array creation
 - ▶ More info:
https://docs.nvidia.com/vpi/architecture.html#arch_array



Array of Harris Corners' (x,y) keypoints

ALGORITHMS

Overview

Algorithms

Gaussian Filter
Stereo Disparity
FFT
...

- ▶ Represent computations on buffers
- ▶ Same algorithm implemented in several backends
- ▶ Different implementations are functionally equivalent
- ▶ Easy to switch implementation among available backends

Algorithm	Backend Support			
	CPU	CUDA	PVA ¹	VIC ²
Box Filter	yes	yes	yes	no
Bilateral Filter	yes	yes	no	no
Gaussian Filter	yes	yes	yes	no
Gaussian Pyramid Generator	yes	yes	yes	no
Convolution	yes	yes	yes	no
Separable Convolution	yes	yes	yes	no
Convert Image Format	yes	yes	no	yes
Rescale	yes	yes	no	yes
Remap	yes	yes	no	yes ¹
Perspective Warp	yes	yes	no	yes ¹
FFT	yes	yes	no	no
Inverse FFT	yes	yes	no	no
Lens Distortion Correction	yes	yes	no	yes ¹
Stereo Disparity Estimator	yes	yes	yes	no
KLT Feature Tracker	yes	yes	yes	no
Harris Corner Detector	yes	yes	yes	no
Temporal Noise Reduction	no	yes	no	yes
Pyramidal LK Optical Flow	yes	yes	no	no

Algorithms available on vpi 1.0

ALGORITHMS

Basic usage

Algorithms

Gaussian Filter
Stereo Disparity
FFT
...

- ▶ Represent computations on buffers
- ▶ Same algorithm implemented in several backends
- ▶ Different implementations are functionally equivalent
- ▶ Easy to switch implementation among available backends
- ▶ Backend must be enabled in the stream

```
/* Apply a 3x3 Box Filter */  
1. vpiSubmitBoxFilter(stream, VPI_BACKEND_CUDA,  
2. imgInput, imgTemp, 3, 3,  
3. VPI_BORDER_CLAMP);  
  
/* Rescale the filtered image */  
4. vpiSubmitRescale(stream, VPI_BACKEND_VIC,  
5. imgTemp, imgResult,  
6. VPI_INTERP_LINEAR,  
7. VPI_BORDER_CLAMP, 0);
```


ALGORITHMS

Algorithm payload

Algorithms

Gaussian Filter
Stereo Disparity
FFT
...

- ▶ Some algorithms need auxiliary resources
- ▶ Such resources are encapsulated in a `VPIPayload` object
- ▶ User is responsible for managing payload lifetime
- ▶ Same payload can't be used concurrently

```
    /* Create stereo payload */
1.  VPIPayload stereo;
2.  vpiCreateStereoDisparityEstimator(VPI_BACKEND_PVA,
3.                                     480, 270,
4.                                     VPI_IMAGE_FORMAT_U16,
5.                                     NULL, &stereo);

    /* Configure some algorithm parameters */
6.  VPIStereoDisparityEstimatorParams params;
7.  params.windowSize    = 5;
8.  Params.maxDisparity = 64;

9.  /* Submits algorithm to process stereo pair */
10. vpiSubmitStereoDisparityEstimator(stream,
11.                                     VPI_BACKEND_PVA,
12.                                     stereo,
13.                                     imgLeft, imgRight,
14.                                     outDisparity, NULL,
15.                                     &params);

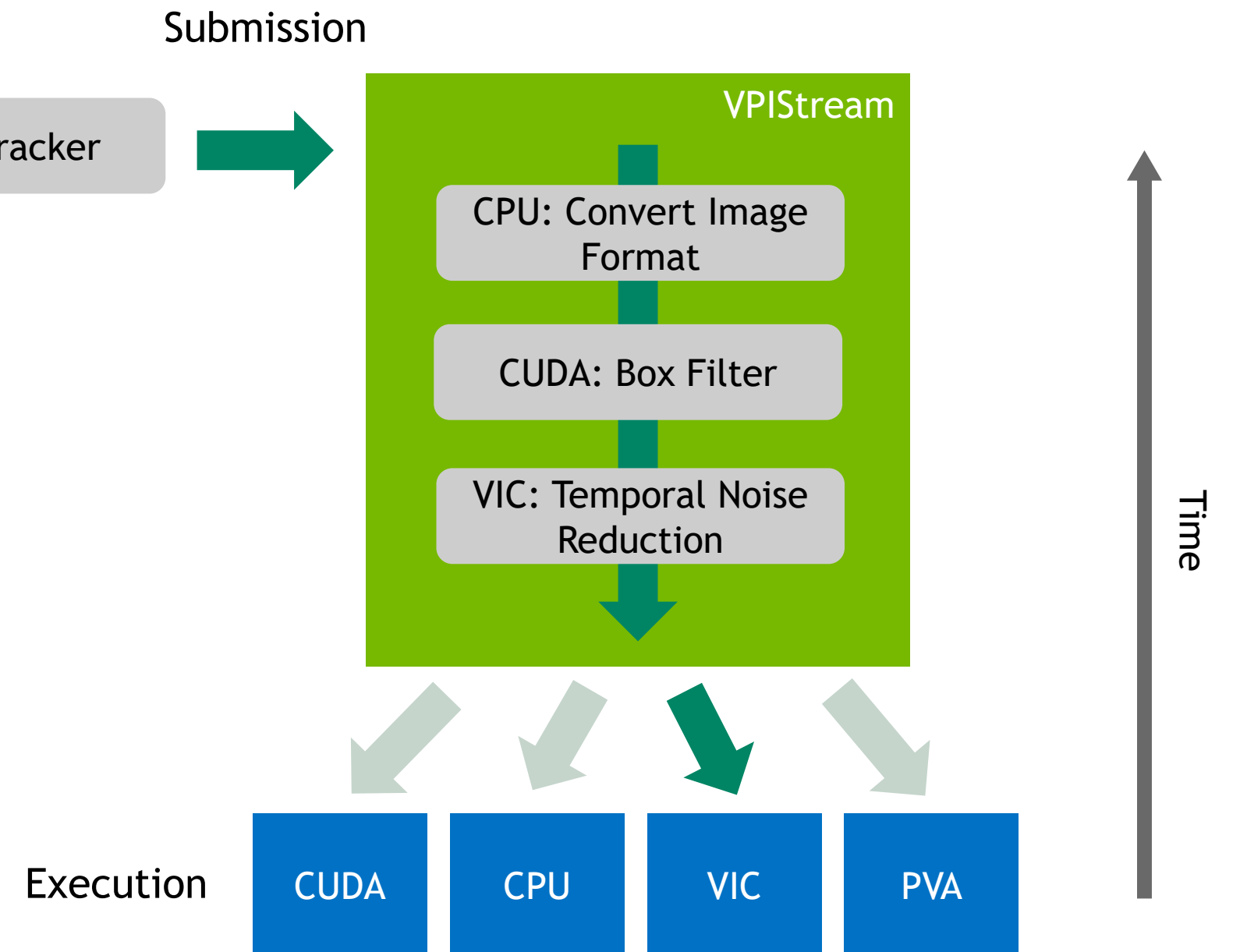
    /* When not needed anymore, destroy payload */
16. vpiPayloadDestroy(stereo);
```

STREAMS

Overview

Streams Execution Queues

- ▶ Represents a task execution queue.
- ▶ Dispatches algorithms to backends for execution.
- ▶ Asynchronous task submission with respect to submission (calling) thread.
- ▶ Use of multiple streams allows for parallel task processing.
- ▶ Destroying a stream implicitly waits for queued tasks to complete.



STREAMS

VPIStream usage

Streams Execution Queues

- ▶ Represents a task execution queue.
- ▶ Dispatches algorithms to backends for execution.
- ▶ Asynchronous task submission with respect to submission (calling) thread.
- ▶ Use of multiple streams allows for parallel task processing.
- ▶ Destroying a stream implicitly waits for queued tasks to complete.

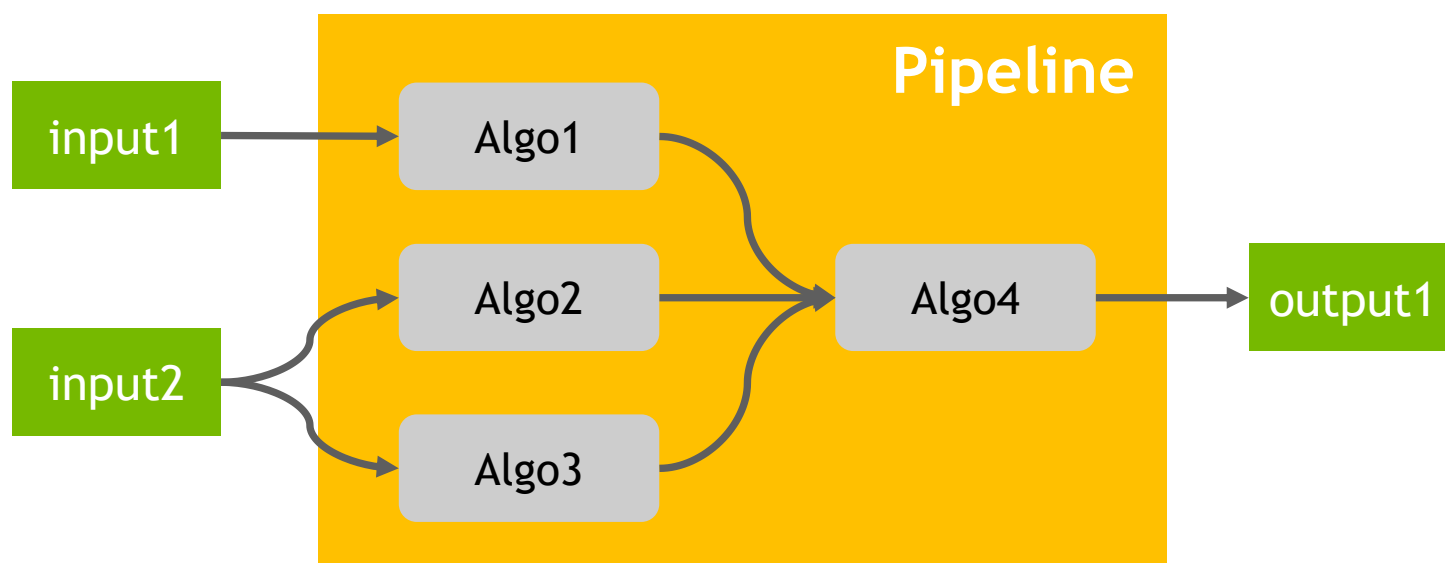
```
    /* Create a stream for algorithm execution in  
       CUDA, PVA and VIC backends. */  
1. VPIStream stream;  
2. vpiStreamCreate(VPI_BACKEND_CUDA | VPI_BACKEND_PVA |  
3.                   VPI_BACKEND_VIC,  
4.                   &stream);  
  
    /* Apply a 3x3 Box Filter */  
5. vpiSubmitBoxFilter(stream, VPI_BACKEND_CUDA,  
6.                   imgInput, imgTemp, 3, 3,  
7.                   VPI_BORDER_CLAMP);  
  
    /* Rescale the filtered image */  
8. vpiSubmitRescale(stream, VPI_BACKEND_VIC,  
9.                   imgTemp, imgResult,  
10.                  VPI_INTERP_LINEAR,  
11.                  VPI_BORDER_CLAMP, 0);  
  
    /* Wait until processing ends */  
12. vpiStreamSync(stream);  
  
    /* Destroy/deallocate the stream */  
13. vpiStreamDestroy(stream);
```


EVENTS

VPIEvent usage

Events Synchronization

- ▶ Synchronization primitives
- ▶ Allow cooperation between streams
- ▶ Modeled after NVIDIA CUDA's `cudaEvent_t`



```
/* Create events */
1. VPIEvent event1, event2;
2. vpiEventCreate(0, &event1);
3. vpiEventCreate(0, &event2);

/* Set up stream 1 */
4. vpiSubmitBoxFilter(stream1, ...); /* Algo1 */
5. vpiEventRecord(event1, stream1);

/* Set up stream 2 */
6. vpiSubmitBoxFilter(stream2, ...); /* Algo2 */
7. vpiEventRecord(event2, stream2);

/* Set up stream 3 */
8. vpiSubmitConvolution(stream3, ...); /* Algo3 */
9. vpiStreamWaitEvent(stream3, event1);
10. vpiStreamWaitEvent(stream3, event2);
11. vpiSubmitKLTFeatureTracker(stream3, ...); /* Algo4 */

/* Wait until processing ends */
12. vpiStreamSync(stream3);

/* Destroy events when not needed anymore */
13. vpiEventDestroy(event1);
14. vpiEventDestroy(event2);
```

ERROR HANDLING

VPIStatus usage

- ▶ Most of VPI functions return a VPIStatus
 - ▶ Operation succeeded: VPI_SUCCESS
 - ▶ Operation failed: VPI_ERROR_*
- ▶ vpiGetLastStatusMessage
 - ▶ Returns detailed information about last error in current thread
- ▶ Asynchronous error reporting
 - ▶ Applies to algorithm execution errors
 - ▶ Reported by events and new algorithm submissions to stream

```
    /* Try to create an image with negative width */
1.  VPIImage img;
2.  VPIStatus status = vpiImageCreate(-1, 256,
                                     VPI_IMAGE_FORMAT_NV12_ER,
3.                                     0, &img);

    /* In case of error, */
4.  if(status != VPI_SUCCESS)
5.  {
6.      assert(status == vpiPeekLastStatus());

7.      char buffer[VPI_MAX_STATUS_MESSAGE_LENGTH];
8.      vpiGetLastStatusMessage(buffer, sizeof(buffer));

9.      printf("Error: %s\n  %s\n", vpiStatusGetName(status),
10.           buffer);

    /* prints:
       VPI_ERROR_INVALID_ARGUMENT
       Width must be strictly positive
    */

11.  exit(1);
12. }
```

PROGRAM STRUCTURE

How VPI-based programs are structured

1. Initialization

- ▶ Allocate resources
 - ▶ Images, Arrays, Pyramids
 - ▶ Streams, Events
 - ▶ Contexts, Payloads
- ▶ Pre-compute static/constant data
- ▶ Define initial processing state

2. Main Processing

- ▶ Fetch input data
- ▶ Execute the processing pipeline
- ▶ Wait for processing to complete
- ▶ Consume result data

3. Cleanup

- ▶ Deallocate resources



**EXAMPLE:
STEREO DISPARITY ESTIMATION**

PROBLEM STATEMENT

What is stereo disparity estimation?

- ▶ Retrieve depth information from stereo pair
- ▶ Stereo pair must be rectified
 - ▶ Rows must be aligned
 - ▶ Might need Lens Distortion Correction
- ▶ Use Semi-Global-Matching technique



Left image

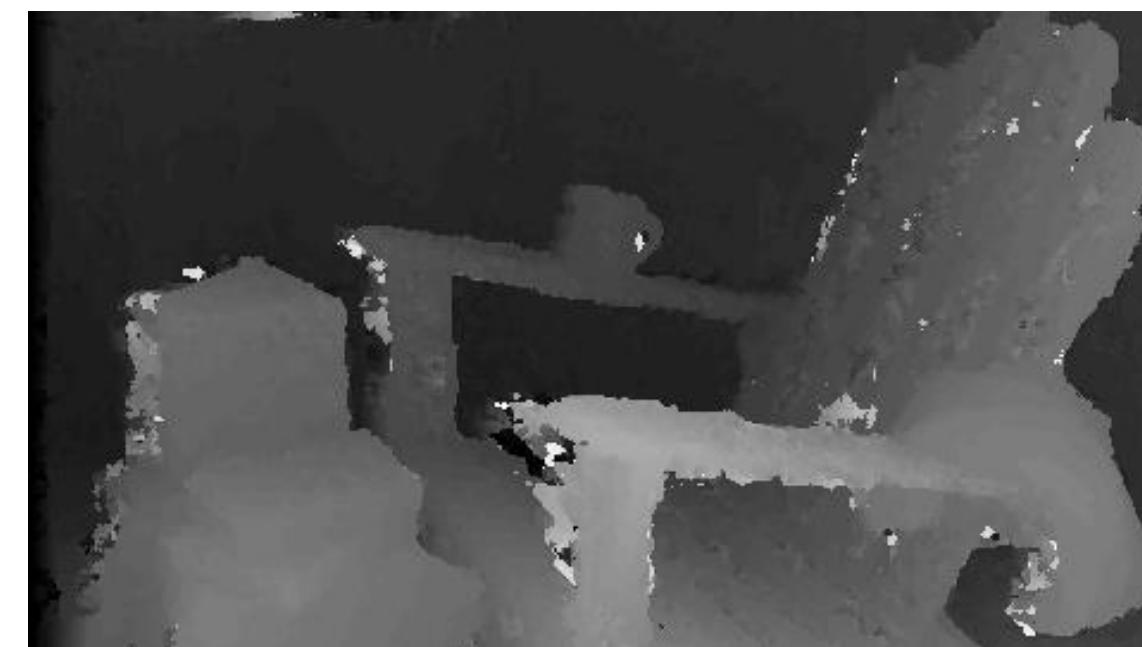


Right image

- ▶ Match pixels on the left image to pixels on the right
- ▶ Calculate similarity of each pixel in the neighborhood, selects most similar.
- ▶ Horizontal pixel displacement = parallax = stereo disparity

▶ Achilles' heel

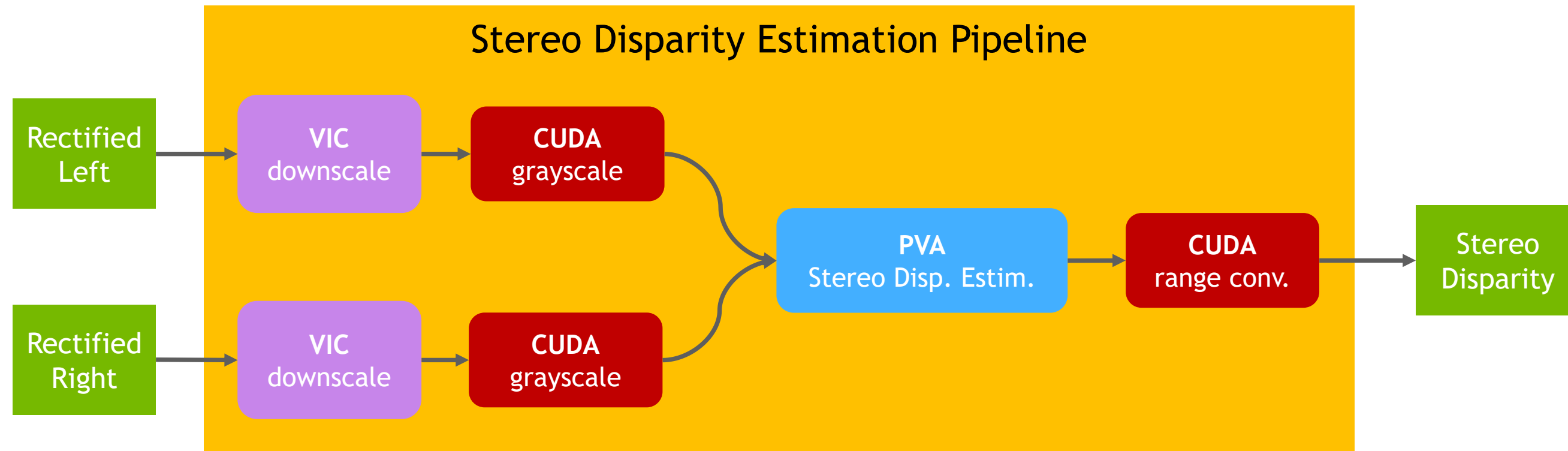
- ▶ Regular patterns/textures
- ▶ Occlusions



Estimated stereo disparity

STEREO PIPELINE

Decomposing the stereo pipeline into stages

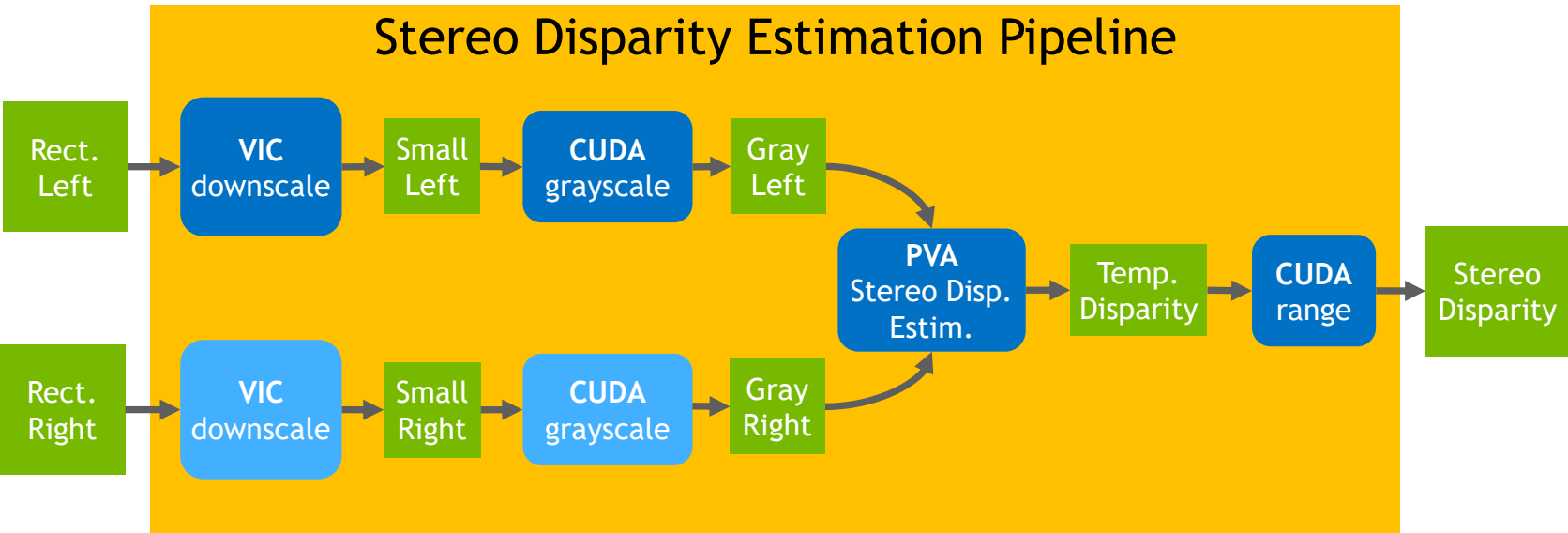
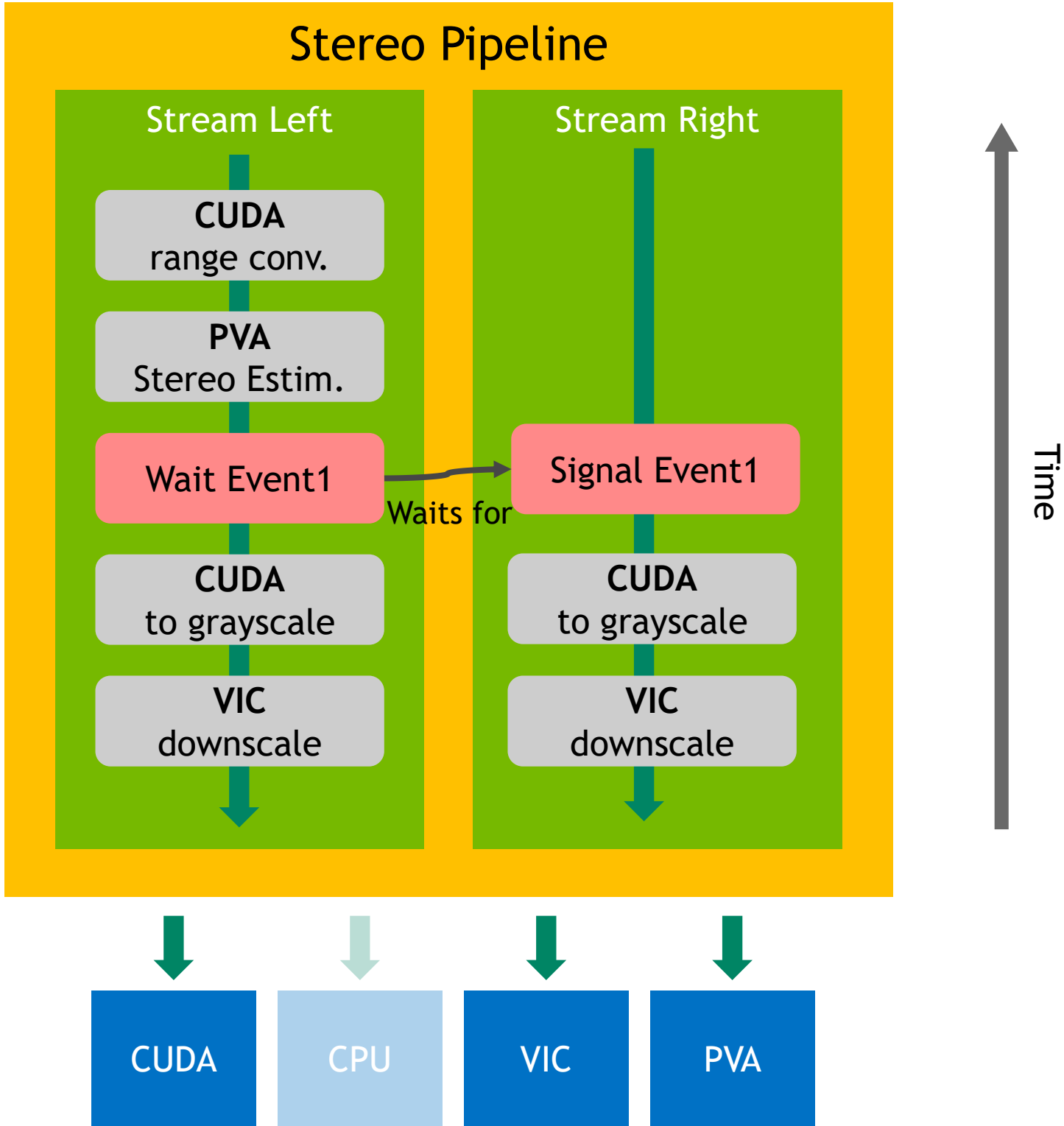


- ▶ As simplification, it expects the stereo pair to be rectified already
- ▶ Stereo Disparity estimation expects its input to be 480x270, U16 grayscale
- ▶ Disparity values' range are rescaled to fit in 0-255 U8 for display
- ▶ Different backends work together towards the same goal

STEREO PIPELINE

Decomposing the stereo pipeline into streams

- ▶ Stereo pair pre-processing defined as parallel stages
- ▶ Reuse left stream for 2nd half of the pipeline
- ▶ Resources needed:
 - ▶ 2x VPIStream
 - ▶ 1x VPIEvent
 - ▶ 8x VPIImage: 2x input, 1x output, 5x intermediate
 - ▶ 1x VPIPayload for PVA stereo processing



INITIALIZATION PHASE

Allocate input VPIImages

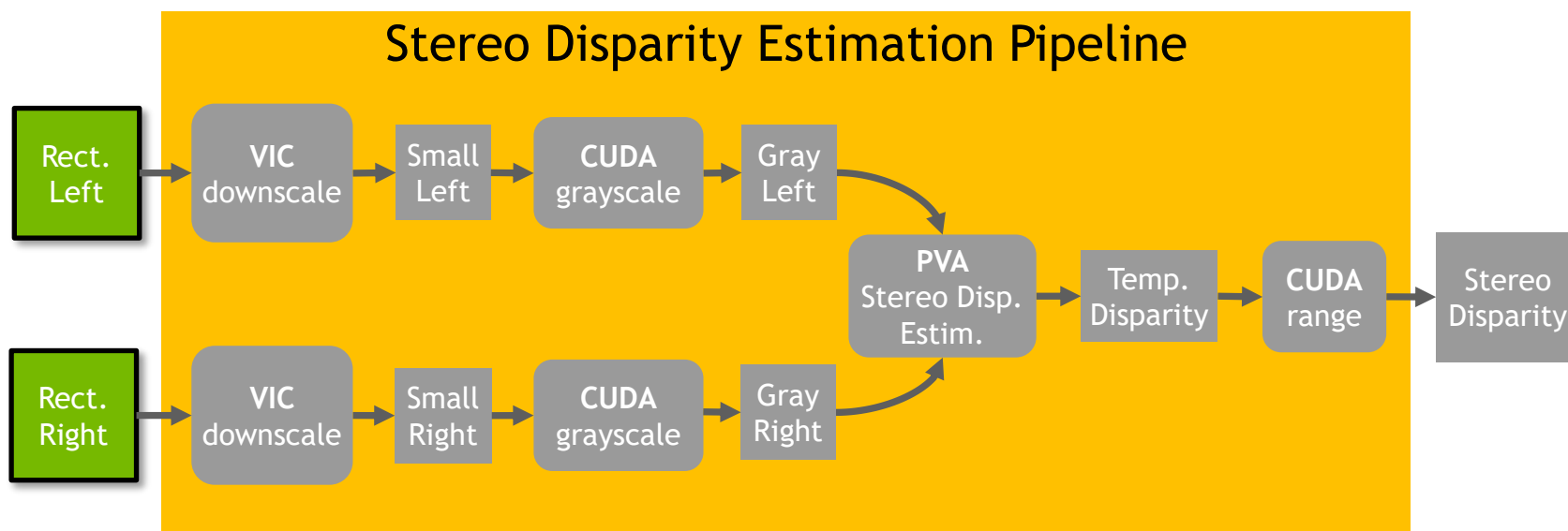
- ▶ Stereo pair comes from 2x cv::VideoCapture
- ▶ Assuming input format is accepted by VIC downscaler
- ▶ Use OpenCV interoperability functions to wrap input cv::Mat into VPIImage
- ▶ Input VPIImage object will be re-used to wrap further input stereo pairs in main loop
- ▶ Error handling omitted for brevity

```
/* Include needed headers */
1. #include <vpi/Image.h>
2. #include <vpi/Stream.h>
3. #include <vpi/Event.h>
4. #include <vpi/OpenCVInterop.hpp>
5. #include <vpi/CUDAInterop.hpp>
6. #include <vpi/algo/Rescale.h>
7. #include <vpi/algo/ConvertImageFormat.h>
8. #include <vpi/algo/StereoDisparity.h>

/* Fetch first stereo pair using
   OpenCV's cv::VideoCapture */
9. cv::Mat cvLeft, cvRight;
10. capLeft.read(cvLeft);
11. capRight.read(cvRight);

/* Wrap them for VPI usage */
12. VPIImage inLeft, inRight;

13. vpiImageCreateOpenCVMatWrapper(cvLeft, VPI_BACKEND_VIC,
14.                                &inLeft);
15. vpiImageCreateOpenCVMatWrapper(cvRight, VPI_BACKEND_VIC,
16.                                &inRight);
```



INITIALIZATION PHASE

Allocate temporary VPIImages

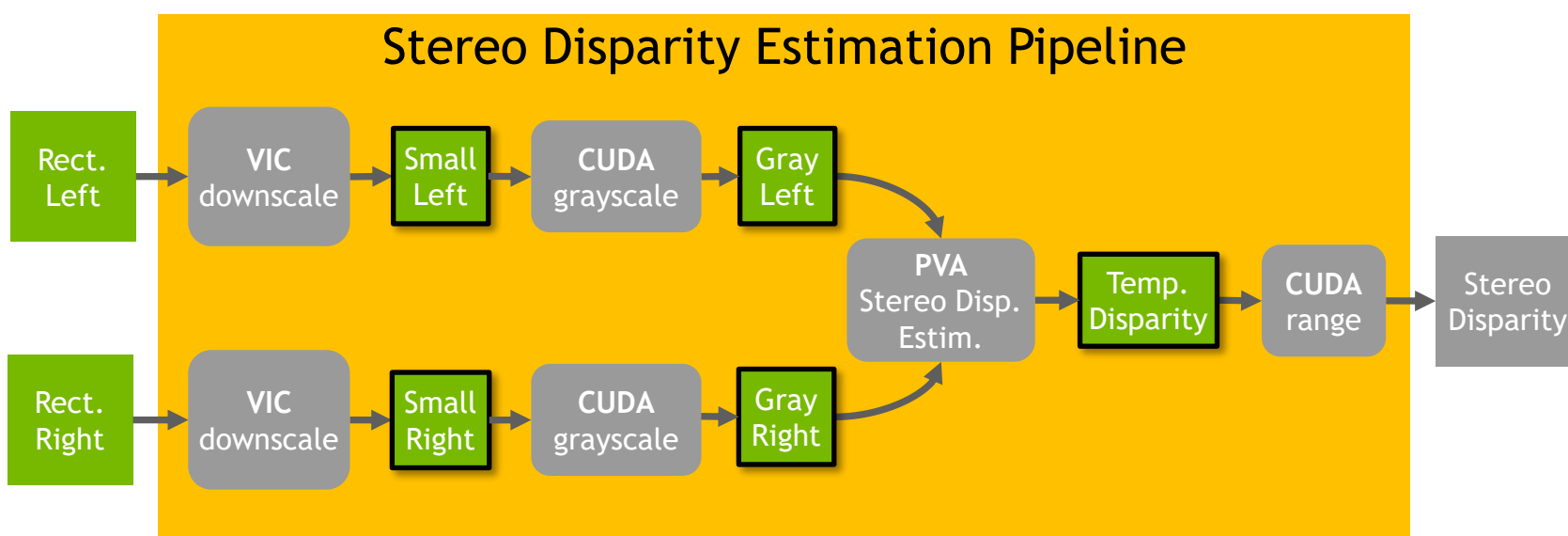
- ▶ Intermediate images must satisfy PVA's Stereo Disparity Estimation restrictions
 - ▶ Input 480x270 U16 (16-bit unsigned values)
 - ▶ Output 480x270 Q11.5 - fixed point with 5 bits for fraction = $2^5 = 32$ equal unit subdivisions

```
    /* Get input image format */
1.  VPIImageFormat inFormat;
2.  vpiImageGetFormat(inLeft, &inFormat);

    /* Allocate downscaled input images */
3.  VPIImage smallLeft, smallRight;
4.  vpiImageCreate(480, 270, inFormat,
5.                VPI_BACKEND_VIC|VPI_BACKEND_CUDA,
6.                &smallLeft);
7.  vpiImageCreate(480, 270, inFormat,
8.                VPI_BACKEND_VIC|VPI_BACKEND_CUDA,
9.                &smallRight);

    /* Allocate grayscale input images */
10. VPIImage grayLeft, grayRight;
11. vpiImageCreate(480, 270, VPI_IMAGE_FORMAT_U16,
12.                VPI_BACKEND_PVA|VPI_BACKEND_CUDA,
13.                &grayLeft);
14. vpiImageCreate(480, 270, VPI_IMAGE_FORMAT_U16,
15.                VPI_BACKEND_PVA|VPI_BACKEND_CUDA,
16.                &grayRight);

    /* Allocate disparity output */
17. VPIImage tempDisparity;
18. vpiImageCreate(480, 270, VPI_IMAGE_FORMAT_U16,
19.                VPI_BACKEND_PVA|VPI_BACKEND_CUDA,
20.                &tempDisparity);
```



INITIALIZATION PHASE

Allocate output VPIImage

- ▶ Access OpenGL's pixel buffer object (PBO) through CUDA's OpenGL interoperability functions.
- ▶ PBO initialization omitted for brevity
- ▶ Wrap resulting CUDA memory buffer into a VPIImage

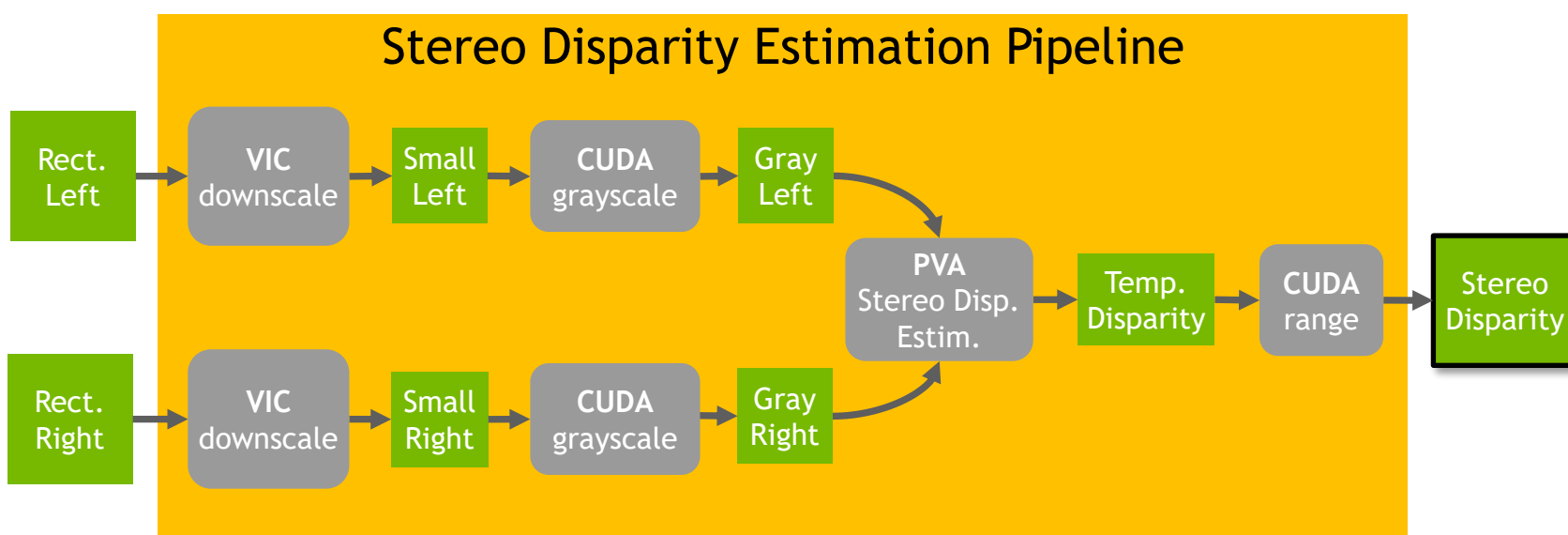
```
/* Create Pixel Buffer Object (PBO) */
1. GLint pbo;
2. glGenBuffers(1, &pbo);
3. /* (Omitted) init `pbo` as a 480x270 8-bit unsigned buffer */

/* Register and map PBO to be accessed by CUDA */
4. glBindBuffer(GL_PIXEL_UNPACK_BUFFER, 0);
5. cudaGraphicsResource_t gres;
6. cudaGraphicsGLRegisterBuffer(&gres, pbo, cudaGraphicsMapFlagsNone);
7. cudaGraphicsMapResources(1, &gres);

/* Retrieve CUDA pointer to PBO memory buffer */
8. void *ptr;
9. size_t len;
10. cudaGraphicsResourceGetMappedPointer(&ptr, &len, gres);

/* Wrap it in a VPIImage */
11. VPIImageData imgData;
12. memset(&imgData, 0, sizeof(imgData));
13. imgData.format = VPI_IMAGE_FORMAT_U8;
14. imgData.numPlanes = 1;
15. imgData.planes[0].width = 480;
16. imgData.planes[0].height = 270;
17. imgData.planes[0].pitchBytes = len/270;
18. imgData.planes[0].data = ptr;

19. VPIImage outDisparity;
20. vpiImageCreateCUDAMemWrapper(&imgData, VPI_BACKEND_CUDA,
21.                               &outDisparity);
```

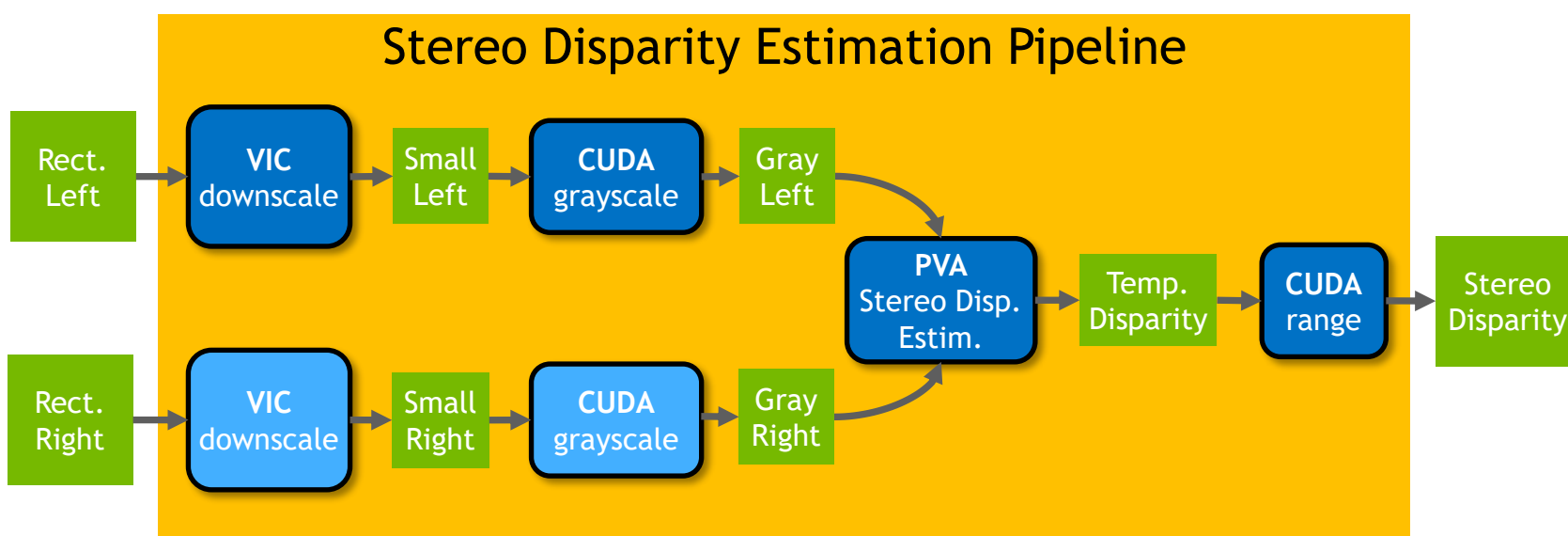


INITIALIZATION PHASE

Allocate VPIStreams, VPIEvent and Stereo VPIPayload

- ▶ Creates Stereo Disparity Estimator payload
 - ▶ Use maxDisparity = 64 for matching search (default)
- ▶ Creates streams with the needed backends enabled
- ▶ Created event object for synchronization
 - ▶ flags = 0, enabled on all backends for simplicity

```
    /* Create stereo payload */  
1.  VPIPayload stereo;  
2.  vpiCreateStereoDisparityEstimator(VPI_BACKEND_PVA,  
3.                                     480, 270,  
4.                                     VPI_IMAGE_FORMAT_U16,  
5.                                     NULL, &stereo);  
  
    /* Create left and right streams */  
6.  VPIStream streamLeft, streamRight;  
7.  vpiStreamCreate(  
8.      VPI_BACKEND_VIC|VPI_BACKEND_CUDA|VPI_BACKEND_PVA,  
9.      &streamLeft);  
  
10. vpiStreamCreate(VPI_BACKEND_VIC|VPI_BACKEND_CUDA,  
                  &streamRight);  
  
    /* Create event to be used in stream synchronization */  
11. VPIEvent evSync;  
12. vpiEventCreate(0, &evSync);
```



MAIN PROCESSING PHASE

Stereo pair pre-processing

- ▶ Retrieve captured stereo pair
- ▶ Wrap them in the existing VPIImages: inLeft/inRight
- ▶ Submit the pre-processing algorithms to the respective streams

```
/* Retrieve captured stereo pair */
1. capLeft.read(cvLeft);
2. capRight.read(cvRight);

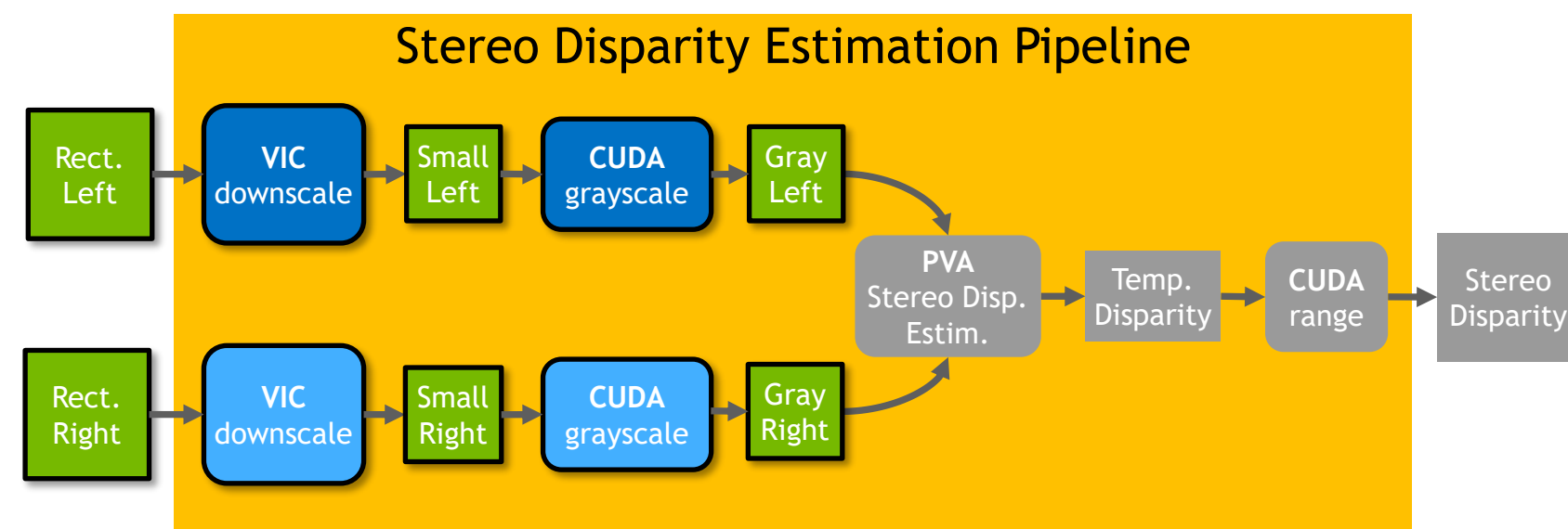
/* Wraps them in the existing inLeft/inRight */
3. vpiImageSetWrappedOpenCVMat(inLeft, cvLeft);
4. vpiImageSetWrappedOpenCVMat(inRight, cvRight);

/* Left image pre-processing */
5. vpiSubmitRescale(streamLeft, VPI_BACKEND_VIC,
6.     inLeft, smallLeft,
7.     VPI_INTERP_LINEAR, VPI_BORDER_CLAMP, 0);

8. vpiSubmitConvertImageFormat(streamLeft, VPI_BACKEND_CUDA,
9.     smallLeft, grayLeft, NULL);

/* Right image pre-processing */
10. vpiSubmitRescale(streamRight, VPI_BACKEND_VIC,
11.     inRight, smallRight,
12.     VPI_INTERP_LINEAR, VPI_BORDER_CLAMP, 0);

13. vpiSubmitConvertImageFormat(streamRight, VPI_BACKEND_CUDA,
14.     smallRight, grayRight, NULL);
```

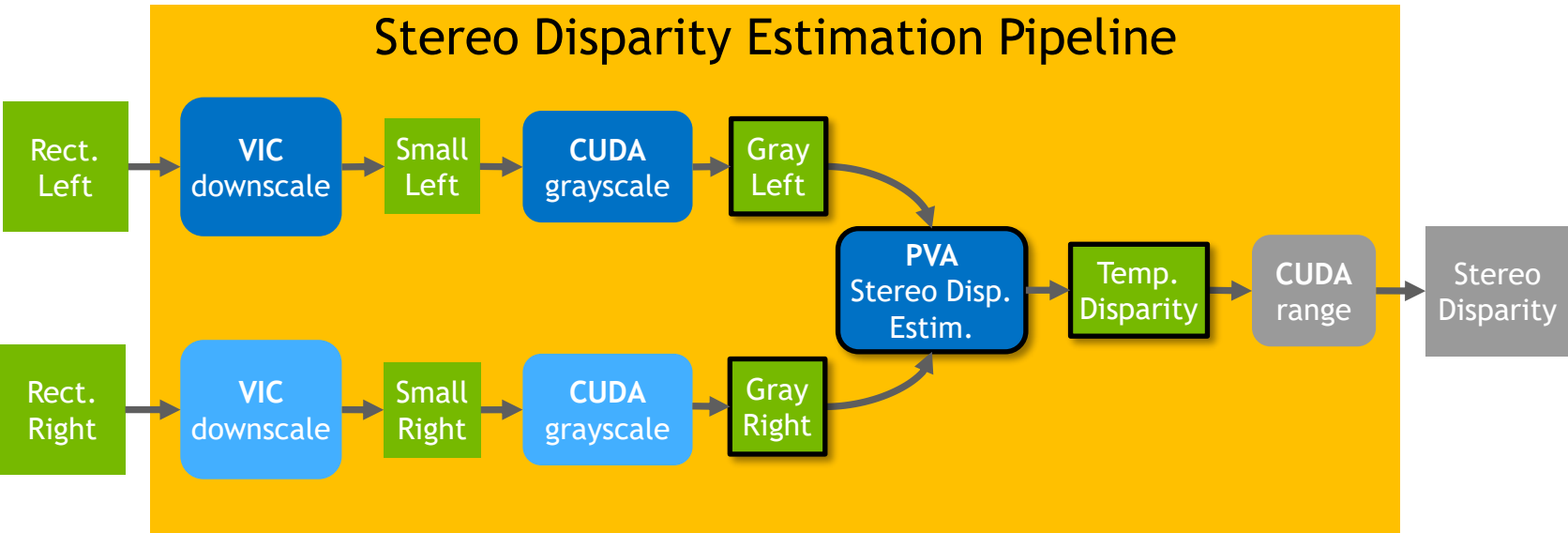


MAIN PROCESSING PHASE

Stereo Disparity Estimation

- ▶ Synchronize streams so that pre-processing tasks are guaranteed to be finished
- ▶ Submit Stereo Disparity Estimation algorithm to stream
 - ▶ Search window size = 5
 - ▶ Maximum disparity for matching search = 64

```
/* Record `streamRight` state on evSync */  
1. vpiEventRecord(evSync, streamRight);  
  
/* `streamLeft` will wait until pre-processing on  
   `streamRight` is finished -> event is signaled */  
2. vpiStreamWaitEvent(streamLeft, evSync);  
  
/* Configure some algorithm parameters */  
3. VPIStereoDisparityEstimatorParams stereoParams;  
4. stereoParams.windowSize = 5;  
5. stereoParams.maxDisparity = 64;  
  
/* Submits algorithm to process stereo pair */  
6. vpiSubmitStereoDisparityEstimator(streamLeft,  
7.                                 VPI_BACKEND_PVA,  
8.                                 stereo,  
9.                                 grayLeft, grayRight,  
10.                                tempDisparity, NULL,  
11.                                &params);
```

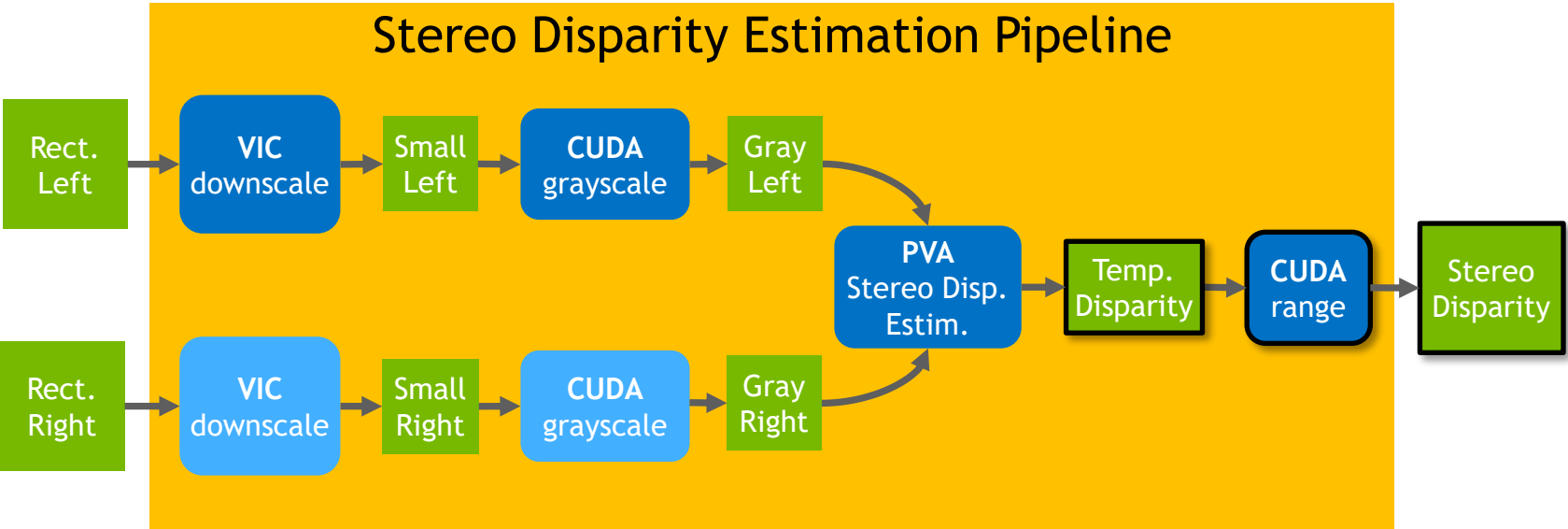


MAIN PROCESSING PHASE

Disparity post-processing

- ▶ Convert disparity from U16 (Q11.5) to U8 using CUDA
- ▶ Wait until processing is finished
 - ▶ Since submit calls are asynchronous, host application can do some other tasks while stream processing isn't finished.
- ▶ Display the resulting PBO (omitted)

```
/* Initialize `cvtParams` with default values */  
1. VPIConvertImageFormatParams cvtParams;  
2. vpiInitConvertImageFormatParams (&cvtParams);  
  
/* Scale disparity range from [0-2047] to [0-255] */  
3. cvtParams.scale = 255.0f/(64*32-1);  
  
/* Submit format conversion from U16 (Q11.5) to U8 */  
4. vpiSubmitConvertImageFormat (streamLeft, VPI_BACKEND_CUDA,  
5.                               tmpDisparity, outDisparity,  
6.                               &cvtParams);  
  
7. /* (Omitted) Host application can do other tasks while  
   stream is processing, as submit calls are  
   asynchronous */  
  
/* Wait until processing is finished */  
8. vpiStreamSync (streamLeft);  
  
9. /* (Omitted) Now PBO can be displayed */
```

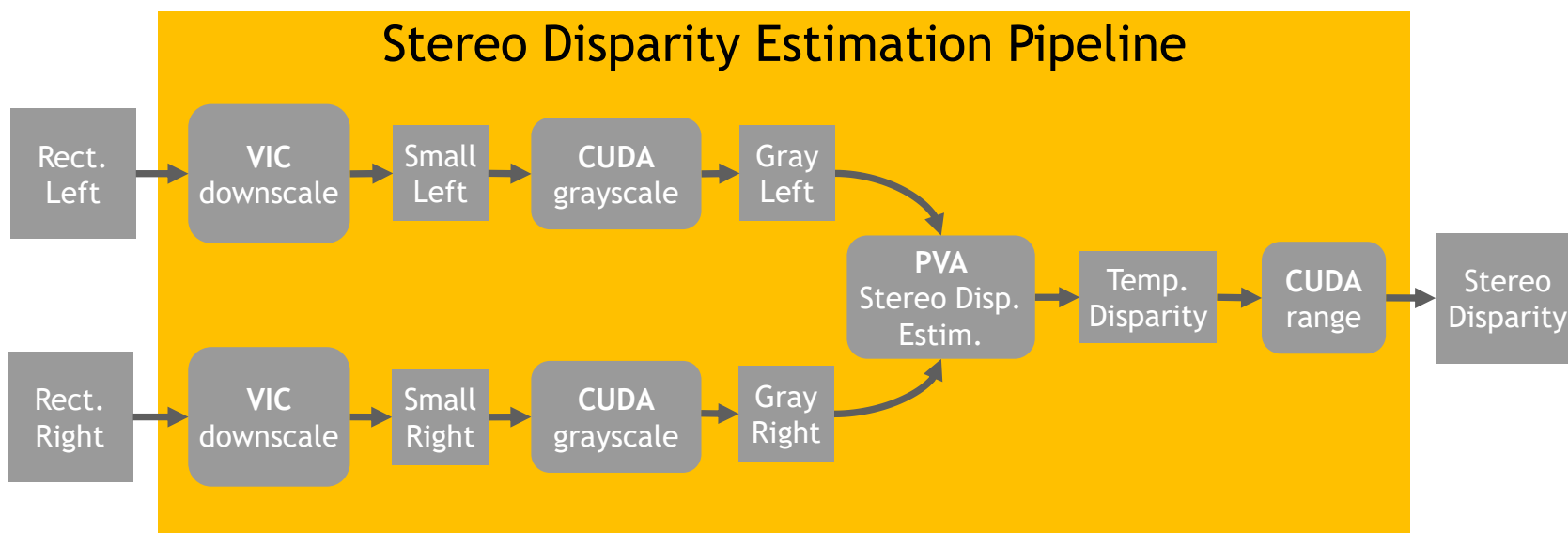


CLEANUP PHASE

Deallocate all resources allocated

- ▶ Prefer to destroy streams first
 - ▶ Destruction implicitly synchronizes it
 - ▶ Guarantees that no tasks are left to be executed
 - ▶ Guarantees that no other objects are being or will be used.
- ▶ Destroy VPIImage wrappers before wrapped buffers

```
    /* Deallocate streams first */  
1. vpiStreamDestroy(leftStream); vpiStreamDestroy(rightStream);  
  
    /* Deallocate event */  
3. vpiEventDestroy(evSync);  
  
    /* Deallocate payload */  
4. vpiPayloadDestroy(stereo);  
  
    /* Deallocate images */  
5. vpiImageDestroy(inLeft);    vpiImageDestroy(inRight);  
6. vpiImageDestroy(smallLeft); vpiImageDestroy(smallRight);  
7. vpiImageDestroy(grayLeft);  vpiImageDestroy(grayRight);  
8. vpiImageDestroy(tempDisparity);  
9. vpiImageDestroy(outDisparity);  
  
    /* Destroy CUDA graphics resource */  
10. cudaGraphicsUnmapResources(1, &gres);  
11. cudaGraphicsUnregisterResource(gres);  
  
    /* Destroy PBO */  
12. glDeleteBuffers(1, &pbo);
```

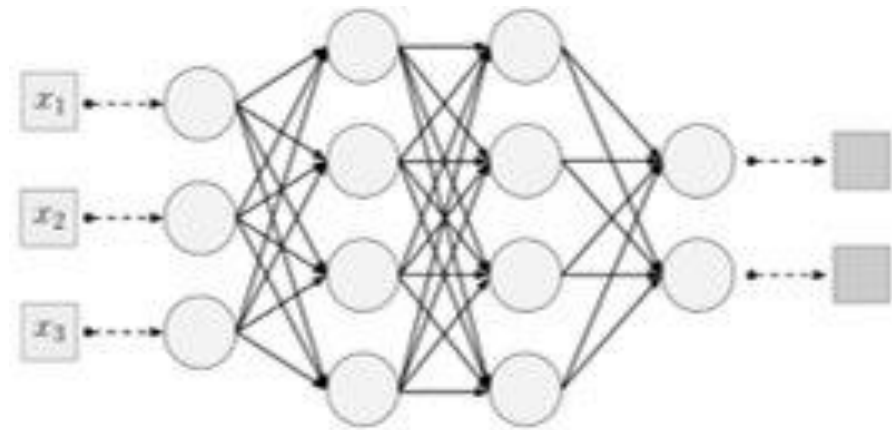




FUTURE WORK

FUTURE WORK

What lies ahead



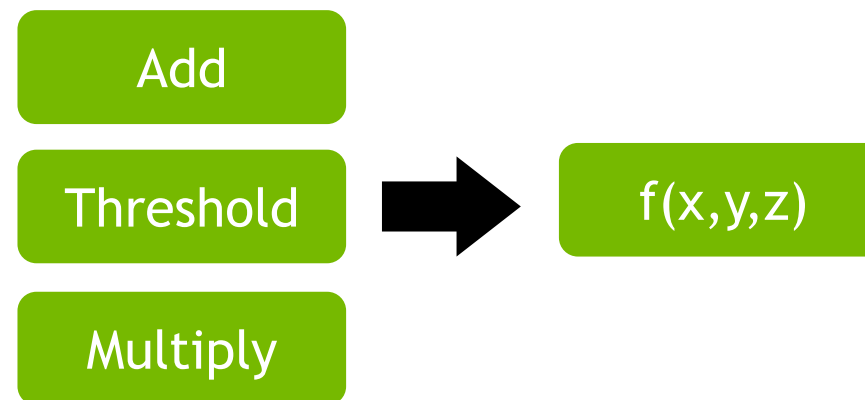
Batching/tensor support for improved DL pre- and post-processing capabilities.



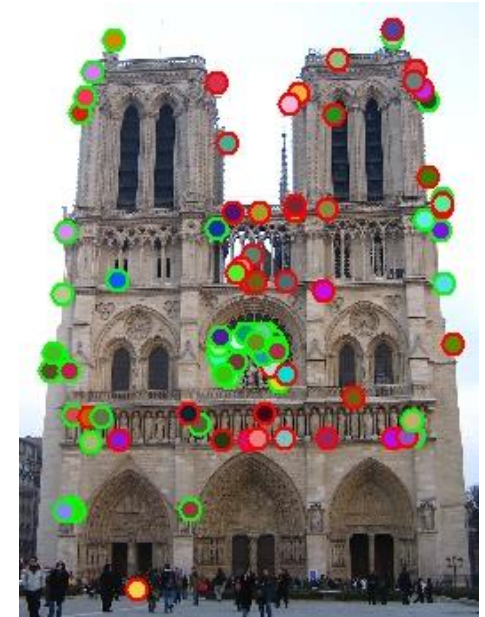
High-level C++ API.



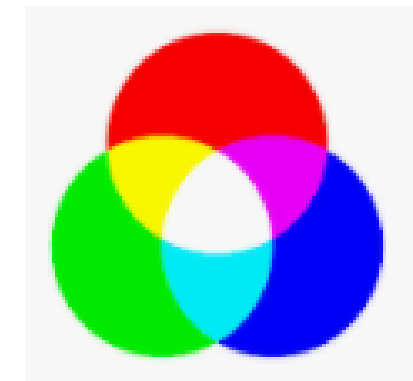
Python bindings.



Algorithm fusion, allowing efficient execution of arbitrary pointwise operations, among others.



More feature detectors, extractors and matchers.



Accurate end-to-end color handling.



Q & A

Thank You!

